# NOAA R&DHPCS ESRL - Boulder

## HPC

**tJet**
Cores: 10128  Nodes: 844
FLOPS: 108.05 TF

**uJet**
Cores: 6048  Nodes: 504
FLOPS: 64.52 TF

**sJet**
Cores: 5440 Nodes: 340
FLOPS: 113.15 TF

**vJet**
Cores: 4608 Nodes: 288
FLOPS: 95.65TF

**xJet**
Cores: 19584 Nodes: 816
FLOPS: 720.69 TF

## Infrastructure

Services:
CAC and RSA Bastion Services
Data Transfer Nodes
NTP
DNS
System Logging
LDAP
Puppet
SNMP
Nessus
TripWire
PerfSonar
MOAB
SSLVPN
ARCSight

## Home File System & Ancillary

JET NAS – NetApp FAS8020 Filers
2 Disk Shelves with 48x900 GB SAS drives
Usable Space: 26.48 TB
Nine File Systems: home (2TB), testhome (2TB), apps
(1TB), test apps (1TB), contrib (130GB), torque (200 GB),
torquelog (100 GB), moab (300 GB), and moablog (100
GB)

Ancillary Systems:
HPC Management
Backup – Amanda
Backup Tape Library: Oracle SL!50 with 2 LTO6 tape drives
Central Syslog
Network/System Monitoring – Zabbix/Nagios

## HPC File Systems

/pan2
File System Type: Panasas
Size: 653 TB

/fs2
File System Type: Lustre
Size: 1.067 PB

/fs3
File System Type: Lustre
Size: 3.1 PB

## Network

WAN:
2x10 GbE connection to NWAVE
2x10 GbE connection to BNOC
Juniper MX80 Border Router
Juniper SRX5600 Firewall

LAN/HPC
2x10 GbE connections WAN/Infrastructure
Juniper EX8202 Core switch
Force10/Dell S50/55/60 cluster core switches
Various vendor edge switches
Production System Connections:
10 GbE
1 GbE
10/100/1000 TX Ethernet

# Detailed HPCS Architecture Boulder

16 January 2017

Boulder Facility Organization and Subsystems Overview

The Boulder HPC systems are divided amongst two rooms in the David Skaggs Research Center (DSRC) in Boulder, CO. Room 2B407 is primarily for low power density air cooled equipment, such as parallel file systems, Front Ends/Login nodes, Batch nodes, Infrastructure servers, and Home File System. There is also a 288 node compute cluster located in 2B407. However, this cluster requires a special air curtain for cold air isolation.

Room GA405 is for high density air cooled equipment such as compute racks. Each compute rack is comprised of 60-70 nodes, IB switches and Ethernet switches. Four of the five compute clusters are located on the GA405 room. Both rooms are connected via 18 active, 150 meter, QDR IB links. So all compute systems reside on the same IB fabric.

Jet is comprised of 5 separate clusters with a total core count of 45,808 cores. A single job cannot span multiple clusters. All clusters can be accessed through the Batch system from the common Front-ends/Login nodes.

| System | Intel Processor | Install Date | Clock (GHz) | Nodes | Total Cores | Memory Per Node (GB) | Memory Per Core (GB) |
|--------|-----------------|--------------|-------------|-------|-------------|----------------------|----------------------|
| tJet | Westmere | 8/2010 | 2.66 | 758 | 9096 | 24 | 2 |
| uJet | Westmere | 11/2011 | 2.66 | 590 | 7080 | 24 | 2 |
| sJet | Sandy Bridge | 8/2012 | 2.6 | 340 | 5440 | 32 | 2 |
| vJet | Ivy Bridge | 8/2014 | 2.6 | 288 | 4608 | 64 | 4 |
| xJet | Haswell | 9/2015 8/2016* | 2.3 | 816 | 19584 | 64 | 2.67 |

* xJet expansion

Jet has approximately 4.82PB of scratch storage. This is across 3 parallel file-systems (2 DDN Lustre and 1 Panasas). Combined peak performance of all three scratch file systems is ~46GB/s. All clusters reside on the same High Speed Interconnect/Infiniband fabric, making the parallel file systems globally accessible from all nodes. There is also a Home File System (HFS) which is globally accessible from all nodes over the Ethernet network.

The Home File System (HFS) and other NFS file systems are being served by a NetApp FAS8120. The HFS is globally accessible from all nodes over the Ethernet network.

The HFS and HPC infrastructure servers are backed up on a daily and weekly basis to a combination of spinning disk and tape. The Open Source software Amanda is used for all backups.

There is a tape storage archive located at the NESCC site in Fairmont, WV. Data can be transferred to and from this archive from Jet's Front Ends. The retention period for the archive is 1-5 years or permanent. All projects with a compute allocation are allowed to store data on the archive.

Data transfer services are provided by a combination of 10GbE connected Data Transfer Nodes (DTN's) and the 10GbE connected Front Ends. The DTN's are used when a data transfer is initiated from outside of Jet. The Jet FE's are used when a data transfer is initiated from Jet.

Jet's Batch system is comprised of three separate pieces of software. Torque is the resource manager and is responsible for job submission, node health status and job launching. Moab is the scheduler and is responsible for scheduling jobs based on priority and managing dedicated system resources with reservations. MAM in the allocation manager and is responsible for managing the monthly allocation of core hours for each project. MAM and Moab are also used for the daily and monthly reporting of system utilization.

# Diagrams

NOAA/Boulder

## 2B201

2B201
~08/01/2016

HFIP p8

**N** ↑

File: Rooms_HFIPp8_2016.current.vsd

Column labels (left to right): AA AB AC AD AE AF AG AH AI AJ AK AL AM AN AO AP AQ AR AS AT AU AV AW AX AY AZ BA BB BC BD BE

Row labels (top to bottom): 01 02 03 04 05 06 07 08 09 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30

- EFSL
- NGSD
- CRAC 8
- T
- UB2A-3
- EFSL-2
- UB2A-1
- T
- UPS UB2 A
- RAMP
- Und. floor Ramp
- NIST T950
- NIST net rack
- nJet TBD
- vJet R85
- vJet R84
- vJet R83
- vJet Spine R82
- vJet R81
- vJet R80
- CRAC 9
- Air Curtain
- Air Curtain
- Col  Area
- Old HSM MD
- OLD DDN 9550
- OLD BDN disk
- Panfs /pan2
- Panfs /pan2
- UB2A-5
- UB2A-2
- Cold  Isle
- SL-150
- ESRL Tape
- ESRL Tape
- ESRL Tape
- Empty Rack
- ESRL/GSD Area
- 811 810 809 808 807    504 503 502 501
- Hot  Isle
- New Net App
- Old Net App
- Juniper
- Patch
- New Net
- New Infra
- Mgmt Serv
- IB Aggr Mgmt
- 610 609 608 607   Patch Panel   603 602 601
- Cold  Isle
- UPS UB2 B
- MIC Test rack
- Empty Rack
- Empty Rack
- /lfstmp Lustre
- nJet Test bed
- UB2B-2   710   708 707   704 703 702 701
- T
- UB2B-1,2
- CRAC 7
- Hot  Isle
- 811 810 809 808 807 806 805 804
- /lfs3 Lustre
- /lfs3 Lustre
- /lfs2 Lustre
- /lfs2 Lustre
- RAMP
- N ↑
- T
- UPS UB2 C
- UB2C-1,2
- UB2A-4
- CRAC 10

NOAA/Boulder

# GA405
## ~08/01/2016

HFIP p8

N

- = xJet, HFIP p8
- = tJet, HFIP p3b
- = uJet, HFIP p4
- = sJet, HFIP p5
- = xJet, HFIP p7

= 2ft x 2ft Floor Tile

Landing

UAG A
UPS-2 Batt | UPS-2 Batt | UPS-2 Batt
Cold Air Deflector
LAG-B  CA
CRAC-17
UAG-A
CRAC-18
UAG-B
CRAC-21

UPS #2
UAG 2
500kVA

CRAC-16

UAG-E

ADA Ramp

UAG1

Fan Coil A/C
Transformer
Work/Storage AREA
Transformer 2x
Transformer 2x

Clean Agent

Bldg. Contour Wand

UAG-D

XDP 1 1000's
XDP 2 Even's

CRAC-19

UPS-1 Batt

UAG-C

Clean Agent

CRAC-20

UPS #1
UAG 1
500 kVA

UPS-1 Batt
UPS-1 Batt
UPS-1 Batt

CA

F10 N-IDF

4128 #1 tJet R40 | 4127 #2 tJet R41 | 4126 #3 tJet R42 | 4125 #4 tJet R43 | 4124 #5 tJet R44 | 4123 #6 Spine R45 | 4122 #7 tJet R46 | 4121 #8 tJet R47 | 4120 #9 tJet R48 | 4119 #10 tJet R49 | 4118 #11 tJet R50

4117 xJet exp R39 | 4116 xJet exp R38 | 4115 xJet exp R37 | 4114 xJet Spine exp R36 | 4113 xJet exp R35 | 4112 xJet exp R34 | 4111 xJet exp R33 | 4110 sJet R76 | 4109 sJet R75 | 4108 sJet R74 | 4107 sJet Spine AGG R73 | 4106 sJet R72 | 4105 sJet R71 | 4104 sJet R70 | 4103 uJet R62 | 4102 uJet R61 | 4101 uJet R60

4217 xJet R90 | 4216 xJet R91 | 4214 xJet Spine R92 | 4213 xJet exp R93 | 4212 xJet exp R94 | 4211 xJet exp R95 | 4210 xJet exp R96 | 4206 uJet R68 | 4205 uJet R67 | 4204 uJet R66 | 4203 uJet R65 | 4202 uJet Spine R64 | 4201 uJet R63 | 4200 uJet old t-13 R52

4311 #12 tJet R51

4301 Tgpu uJet R53

File: Rooms_HFIPp8_2016.current.vsd

Page: GA405

# Boulder – External Connectivity



**Legend**

| | |
|---|---|
| Ten Gigabit Ethernet | ———— |
| Gigabit Ethernet | ———— |
| 10/100/1000TX Ethernet | ———— |

# Boulder – Infrastructure
## RM 2B201

140.208.160.128/26
(DMZ7)

GSD NAS

HPCS
Compute/File
Systems

Multiple Connections

Juniper
EX8208
bldrcore

x2

0500
Boundary

Boulder
External Connections

Internet

BNOC

n-wave

HPCS FW
Juniper
SRX5600

BLDR CE
Juniper
MX80

140.208.160.0/26   (DMZ5)

| | | |
|---|---|---|
| **FE** | fe[1-8], tfe[1-2] | |
| **jetadm** | jetadm[3-5] | |
| **jetadm** | jetadm[1-2] | |
| **mgmt** | wms[1-4] | |
| **mgmt** | wms[5-16] | |

sw17-3
jetcore1
sw-mgmt-1

140.208.169.0/25 (903)
140.208.169.0/25 (903)

| DMZ-3 | DMZ-4 | |
|---|---|---|
| **ips/idp** | B-021 | 160 |
| **amanda** | B-014 | 154 |
| **snmp, nessus** | B-012 | 152 |
| **loghost** | B-013 | 153 |
| **qradar-collector** | B-020 | 159 |

Outside
140.208.170.16/28 (899)

.157 | **Perfsonar** | B-018

140.208.168.128/25 (902)

| DMZ-2 | DMZ-4 | |
|---|---|---|
| **idns1, ntp1, ildap1 monitor** | B-003 | 145 |
| **idns2, ntp2, ildap2, cvs, relay2** | B-004 | 146 |
| **Moab3** | B-007 | 149 |
| **Myproxy1** | B-010 | 150 |
| **Myproxy2** | B-011 | 151 |
| | B-017 | 156 |

Cisco 2900

DMZ-4
140.208.169.128/25 (904)

| | | |
|---|---|---|
| .130 | **Puppet** | B-015 |
| | **Avocent KVM** | ipkvm2 |
| | **Mon/Key** | ipkvm3 |
| .161 | **Net Sec Mgr** | B-022 |
| | **PDU 1-2** | |
| .141 | **APC Transf'm SW** | transw-001 |
| .142 | **APC Transf'm SW** | transw-002 |

140.208.170.32/28 (899)

140.208.168.0/25 (902)

| DMZ-8 | DMZ-4 | |
|---|---|---|
| **rsa-02** | rsa-02 | 168 |
| **rsa-03** | rsa-03 | 169 |

x2

x2

140.208.168.0/25 (902)

140.208.168.0/25   (DMZ1)
140.208.170.48/28 (DMZ14)

RDHPCS
SSL-VPN1

RDHPCS
SSL-VPN2

| DMZ-1 | DMZ-4 | |
|---|---|---|
| **Bastion** | B-001 (B-016 HW) | 143 |
| **edns1,eldap,relay1,www** | | |
| **helpdesk,docs** | B-002 | 144 |

| DMZ-1 | | |
|---|---|---|
| **jetscp1** | jetscp1 | |
| **jetscp2** | jetscp2 | |
| **testscp1** | testscp1 | |
| **testscp2** | testscp2 | |

| NOAA R&DHPCS Logical Network Diagrams | |
|---|---|
| Date: 8/30/16 | Version 5.8 |
| Drawn For: National Oceanic and Atmospheric Administration | |
| | Pg: 4 |

Page: Boulder Infrastructur

# Boulder – HPCS
# RM GA405

IB Aggr 0,1,2 ──3x── Netgear48p sw28-2

sw-mgmt-1
10.178.17.0/24
10.178.28.0/24
140.208.160.0/26

**tJet**

DRAC

wms9
wms10

10.178.45.0/24
140.208.160.0/26

8 racks

490 Compute
IB Edge/Leaf (23)
IB Spines (18)

14 switches
2 in R40-44, 46-47
10.178.$rack.0/24

sw45-1
10.178.45.0/24

Dell 48p

Dell S50N
Tcore1
10.178.45.0/24

DRAC

wms7
wms8

10.178.45.0/24
140.208.160.0/26

4 racks

268 Compute
IB Edge/Leaf (13)

8 switches
2 in R48-51
10.178.$rack.0/24

Dell 48p

Dell S50N
Tcore2
10.178.45.0/24

**uJet**

1 rack

22 Compute
IB Edge/Leaf (1)

1 switch in R53
10.178.$rack.0/24

1 rack

64 Compute
IB Edge/Leaf (3)

2 switches in R52
10.178.$rack.0/24

9 racks

504 Compute
IB Edge/Leaf (26)
IB Spines (14)

16 switches
2 in R60-63, 65-68
10.178.$rack.0/24

Dell 48p

sw64-1
10.178.64.0/24

Dell 48p

Dell S50N
Ucore1
10.178.64.0/24

wms11
wms12
10.178.64.0/24
140.208.160.0/26

DRAC

Boulder
HPCS 2B201

Juniper
EX8208
bldrcore

12x (DRAC + External IP)

**sJet**

9 switches
1 in R71,74,76, 2 in R70,72,75
10.178.$rack.0/24

7 racks
340 nodes
34 chassis's(10 nodes per)

Dell S50N
Score1
10.178.73.0/24

Brocade
FCX48G

ISCB

Node 10

Compute
Node 1 | Chassis

IB Edge/Leaf (17)
IB Spines (7)

wms13
wms14
10.178.73.0/24
140.208.160.0/26

**xJet**

13 switches
2 in R90,91,93,94,95,96 1 in R92
10.178.$rack.0/24

7 racks
396 Compute nodes
4 Big Memory nodes
100 chassis's(4 nodes per)

Dell S50N
Xcore1
10.178.92.0/24

Brocade
FCX48G

Node 4

Compute
Node 1 | Chassis

IB Edge/Leaf (17)
IB Spines (6)

10.178.92.0/24
140.208.160.0/26

wms15
wms16

13 switches
2 in R33,34,35,37,38,39, 1 in R36
10.178.$rack.0/24

7 racks
416 Compute nodes
104 chassis's(4 nodes per)

Dell S55
Xcore2
10.178.36.0/24

Brocade
FCX48G

Node 4

Compute
Node 1 | Chassis

IB Edge/Leaf (17)
IB Spines (6)

10.178.36.0/24
140.208.160.0/26

wms5
wms6

## Legend

| | |
|---|---|
| Ten Gigabit Ethernet | |
| Gigabit Ethernet | |
| 10/100/1000TX Ethernet | |

| NOAA R&DHPCS Logical Network Diagrams | |
|---|---|
| Date: 8/30/16 | Version: 5.8 |
| Drawn For: National Oceanic and Atmospheric Administration | |
| | Pg: 5 |

# Boulder – HPCS
# RM 2B201

**nJet Test System**

2 switches
10.178.18.0/24

SMC 48p

62 Compute nodes 1 rack

**Compute**
IB Edge/Leaf (3)
IB Spines(2)
tbqs1
tbqs2
tfe1
tfe2
testscp1
testscp2

Dell S50N Testcore1
10.178.20.0/24
10.178.18.0/24

4x

wms3
wms4
10.178.20.0/24
140.208.160.0/26

4x

8x

**vJet**

Dell S50N
vcore1
10.178.20.0/24

10 switches
2 in R80,81,83,84,85
10.178.$rack.0/24

Brocade FCX48G

2x

288x

6 racks
288 nodes
72 chassis's(4 nodes per)

Node 4
...
**Compute** Node 1 Chassis
IB Edge/Leaf (12)
IB Spines (6) 1 rack

**Jetcore**

wms1
wms2
10.178.126.0/24
140.208.160.0/26
4x (DRAC + External IP)

IB swio-[2-5]

IB Aggr 3,4

lbms3
lbms4
DRAC
10.178.17.0/24

**Dell S50N
Jetcore1**
140.208.160.0/26
10.178.17.0/24
10.178.126.0/24
10.178.20.0/24

2x

8x

2 nodes
10.178.126.0/24
140.208.168.0/25
**Jetscp**

8 nodes
10.178.126.0/24
140.208.160.0/26
**FE**

10.178.17.0/24
**Sherwood**

2x

Netgear 48p sw17-2
10.178.17.0/24

DRAC

**Jetbqs[3-4]**

2 nodes

2x

MD3420-3
MD3420-4
MD3420

**SL150**
Backup Library
Brocade 300 (FC)

10.178.17.0/24
140.208.160.0/26

jetadm[1-2]
2 nodes

MD3420-1
MD3420-2
MD3420

4x

2x

sw17-3 Dell 48p
10.178.17.0/24
140.208.160.0/26

3x DRAC

**Jetadm[3-5]**

Backup

2x Backup

Backup

File: NOAA0500-Network-Diagrams-v5.8.vsd

GSD NAS

Boulder
Infrastructure

2x

**Juniper
EX8208
bldrcore**

2x

2x

8x

2x

2x

2x

2x

3x

26x

2x

2x
2x

4x

Boulder
HPCS GA405

**JetNAS**
10.178.17.0/24

jetnas1
jetnas2
NetApp
FAS8020

Page 4 of 4

BLDR CE MX80

Boulder
External
Connections

HPCS FW SRX5600

**Panasas**

**Dell S50N
sw123-1**
10.178.123.0/24

14 shelves
/pan2
shelves

28x

**Dell S4810**
sw123-2
10.178.123.0/24

2x

8x

14x

8 nodes
/pan2
IB router

14x

**Lustre**

lfs2

Directmon (MGS)

2 nodes
**MDS**

MDT(EF3015)

Dell 24p sw16-4
10.178.16.0/24

28x

8 nodes
**OSS**

OST (SFA12k-40)

Dell 24p sw16-5
10.178.16.0/24
Servers (DRAC + Eth0)

16x

lfs3

2 nodes
**MDS**

MDT(SFA7700)

Dell 24p sw16-6
10.178.16.0/24
SFA Controllers

16x

8 nodes
**OSS**

OST(4 x SFA7700)

Legend

Ten Gigabit Ethernet

Gigabit Ethernet

10/100/1000TX Ethernet

NOAA R&DHPCS
Logical Network Diagrams

Date: 8/30/16 | Version: 5.8

Drawn For: National Oceanic
and Atmospheric Administration

Pg: 6

Page: Boulder HPCS RM 2B201

# IB HSN HL Layout

**HFIP p8**
**08/18/2016**

| GA405 | 2B201 |
|---|---|

QDR Aggr 2
@2 QDR

QDR Aggr 5
@2 QDR

Note:
Uses 9 groups of
2 between rooms

QDR Aggr 1
@2 QDR

QDR Aggr 4
@2 QDR

18 cables, (2 spare)
150m, 72 GB/s

QDR Aggr 0
@2 QDR

QDR Aggr 3
@2 QDR

**tJet** Sub-tree
**uJet** Sub-tree
**sJet** Sub-tree
**xJet** Sub-tree
**nTest** Testbed
**I/O & Service**
**vJet** Sub-tree

18 cables 72 GB/s
14 cables 56 GB/s
14 cables 56 GB/s
24 cables 96 GB/s
6 cables 24 GB/s
27 cables 108 GB/s
12 cables 81 GB/s

Spines — Spines — Spines — Spines — Spines — Edges — Spines

Edges — Edges — Edges — Edges — Edges — Edges

Nodes — Nodes — Nodes — Nodes — Nodes

Subnet Mgrs | Service Nodes | /pan2 | /lfs tmp | /lfs2 | /lfs3

Nodes

| GA405 | 2B201 |
|---|---|

IB, FDR, 4x, 56Gb/s line, 54 Gb/s data  — · — · —
IB, QDR, 4x, 40Gb/s line, 32 Gb/s data  — · · — · · —
FDR IB Switch, 36p
QDR IB Switch, 36p

# IB HSN Layout - 2B201

**HFIP p8**
**08/18/2016**

18 cables + 2 spare
150m   72 GB/S
**To GA405**

Note:
Uses 9 groups of
2 between rooms

Aggr 5
6
18
@2
QDR
12

Aggr 4
6
18
@2
QDR
12

Aggr 3
6
18
@2
QDR
12

**nTest**
**Testbed**
6 cables
24 GB/s

**I/O**
27 cables
108 GB/s

@3
QDR
9
5   iosw 3
22

iosw 4
9
5
22

iosw 5
9
9
18

@1
QDR
Spine 1
Spine 2
3
24

2 spines
6/21=29% BBW
6*3=18 cables

**vJet**
**Sub-Tree**

12 cables
81 GB/s

6 spines
12/24 =50% BBW
6*24=144 cables

2
10   Spine 1 ···· Spine 6
24

@2
FDR

@2
FDR

Edge 1
@3
QDR
Edge 3
6
10
20

9
21

@1
QDR

12
0   Edge 1 ········ Edge 12
24

288 nodes
0 free ports

@1
FDR

Subnet
Mgrs

@1

**Service**
**8 FE**
**2 SCP**

**Service**
**1 RBH**
**4 TST**

**/pan2**
Panasas
8 IO routers

12.5GB/s,
653TB

**/lfstmp**
DDN/
Lustre
2x MDS
2x OSS

**/lfs2**
DDN/Lustre
8x OSS
2x MDS
1xMGS
15 GB/s,
1070TB

**/lfs3**
DDN/Lustre
2x MDS
8x OSS
23 GB/s
3080 TB

62 nodes
28 free ports

Node 1 ···· Node 24   Node 1 ···· Node 24

Node 1   Node 21   Node 1   Node 20

**vJet**
6  spine switches
12 clos switches
24  nodes/clos
288 nodes (408 max)

**nJet testbed**
2  spine switches
3 clos switches
21  nodes/clos
62 nodes

IB, FDR, 4x, 56Gb/s line, 54 Gb/s data  — · — · —
IB, QDR, 4x, 40Gb/s line, 32 Gb/s data  — ·· — ·· —
FDR IB Switch, 36p  (teal hexagon)
QDR IB Switch, 36p  (yellow hexagon)

# IB HSN Layout - GA405

HFIP p8
08/18/2016



**tJet subtree**
- 18 spine switches
- 36 edge switches (42 max)
- 21 or 22 nodes/edge
- 758 tJet nodes
- (882 nodes max @ BBW)

**uJet subtree**
- 14 spine switches
- 28 edge switches (35 max)
- 21 nodes/edge
- 590 uJet nodes
- (735 max @ BBW)

**sJet subtree**
- 7 spine switches
- 17 edge switches
- 20 nodes/clos
- 340 nodes (max @ BBW)

**xJet subtree**
- 12 spine switches
- 34 edge switches (34 max)
- 24 nodes/clos
- 816 nodes (816 max @ BBW)

IB, FDR, 4x, 56Gb/s line, 54 Gb/s data  — · — · —
IB, QDR, 4x, 40Gb/s line, 32 Gb/s data  — · · — · · —
FDR IB Switch, 36p
QDR IB Switch, 36p